

# An Action-based Training Method For *VideoPose3D*

Hao Bai  
ZJU-UIUC Institute  
Haining, Jiaxing

<https://jackgetup.com>

Prof. Gaoang Wang  
ZJU-UIUC Institute  
Haining, Jiaxing

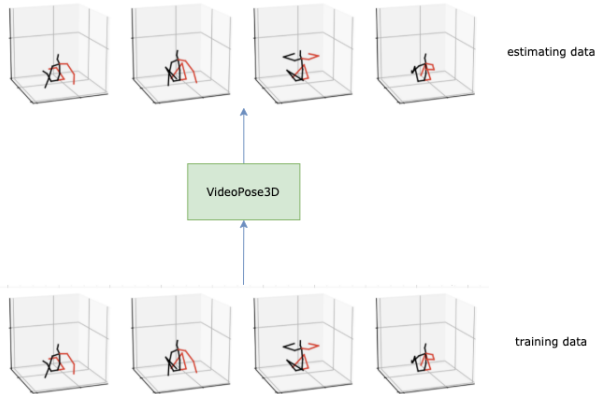
[gaoangwang@intl.zju.edu.cn](mailto:gaoangwang@intl.zju.edu.cn)

## Abstract

*In this essay we perform a homogenized action-based training method based on VideoPose3D. The original paper utilizes a semi-supervised approach on different kinds of subjects and different actions in the network, and we select each action as the only input of the training dataset, and estimate the error in multiple protocols. Our approach is utilized to reduce the error on Human3.6M.*

## 1. Introduction

Vision-based human motion analysis attempts to understand the movements of the human body using computer vision and machine learning techniques. Our work focuses on the improvement of 3D human pose estimation in video. The original *VideoPose3D* model carries out the result shown in Figure 1, and our model changes the idea with Figure 2.



### 1.1. Action-based Learning

Compared to the original model we utilized a different data-processing approach. The original approach mentioned in *VideoPose3D* uses all actions in training data, and the estimating effect was good overall. We utilize the method that we only take one type of action for training, and the total amount for training should be the same. Math-

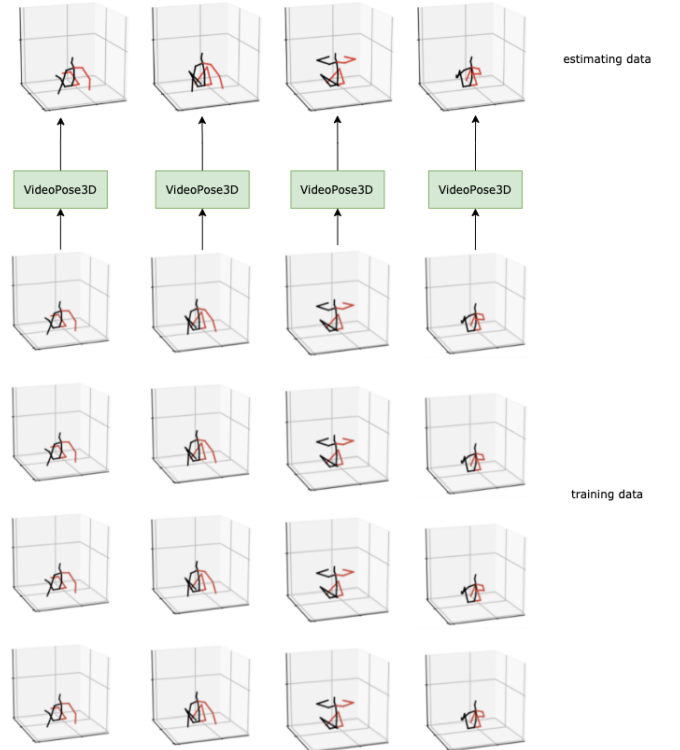
ematically, take the action *Sitting* for example, it can be expressed as

$$\sum_{i \in \text{actions}} f(i) = f(\text{Sitting}) \cdot t \quad (1)$$

where  $f$  stands for the frames of subject, and  $t$  means the number of epochs of our action-based training.

In this approach, we utilize only one time of action, which means we save more data; and we also produce the results in the same time, the evidence can be illustrated by formula (1). After testing, we conclude that our work outperforms the original work by 11%.

The principle of our work is shown below.

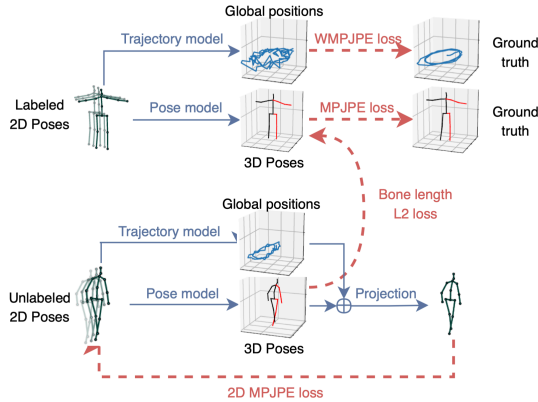


## 2. Related Work

Our work is based on the previous work *VideoPose3D* with an inspiration for its semi-supervised approach. Earlier work has shown the effect of semi-supervised learning as a great way for improving training performance [8].

### 2.1. VideoPose3D

*VideoPose3D* is the current state-of-the-art approach which utilizes a fully convolutional model based on dilated temporal convolutions over 2D keypoints. [4, 5]. In early researches the main solution was to use recurrent neural networks (RNN) [2], but the research group of *VideoPose3D* utilized the convolutional neuron network (CNN) to gain an efficient result on the dataset *Human3.6M*, with the inspiration that many authors had mentioned CNN in temporal models. Some researches has also shown the idea of using a spacial-temporal model being useful [3], but that's not utilized in our research.



### 2.2. Action-based Learning

The idea of homogenized learning was proposed firstly by Action Estimations [7, 6].

### 2.3. Homography Estimation

The idea of using homonegeious actions in a dataset was also inspired from Homography Estimation in some deep images [1].

## 3. Experimental Setup

For the code of this project, please refer to the url below.

[https://github.com/BiEchi/](https://github.com/BiEchi/Pose3dDirectionalTraining)

Pose3dDirectionalTraining.

The detailed experimental setup can also be explored in *README.md* in the github repo. To reproduce the results, you need to test for 15 times for each type of action, and it's preferable to test for multiple times. It's more convenient to use *bash* scripts to save time.

## 4. Results

Shown below are the results of the experiment. Expect a minor error when testing on your own. Note that all unit are millimeter.

### 4.1. Small Number of Epochs

Table 1. MPJPE, 3\*3, unit-epoch

Action	Original, MPJPE	Action-based, MPJPE
Eating	56.34	<b>52.07</b>
Sitting	70.49	70.86

Table 2. Velocity Error, 3\*3, unit-epoch

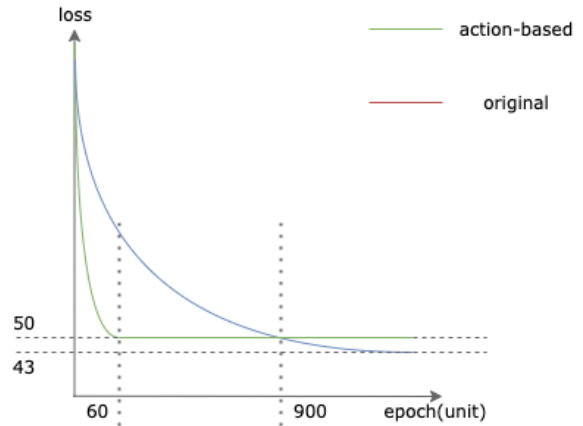
Action	Original, Velo-M	Action-based, Velo-M
Eating	3.14	<b>3.01</b>
Sitting	3.14	<b>2.86</b>

### 4.2. Large Number of Epochs

For a large number of data, we observe that because of the characteristics of the model that *VideoPose3D* utilizes, the epoch refuses to increase after epoch 60 (i.e. epoch 4 of the original training method).

Table 3. Property Cluster, 3\*3, 80-epoch

Type	Original	Eating-based
Time Consumed	81173 sec	56921 sec
MPJPE-Eating	43.30 mm	49.27 mm
Velo-M-Eating	2.24 mm	2.56 mm



## 5. Conclusion

According to the results shown above, we conclude that the action-based model has a better performance on small epochs because they learn faster than the original one, though it has a worse performance on the top performance.

## References

- [1] Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Deep image homography estimation. *arXiv preprint arXiv:1606.03798*, 2016.
- [2] Kyoungoh Lee, Inwoong Lee, and Sanghoon Lee. Propagating lstm: 3d pose estimation based on joint interdependency. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 119–135, 2018.
- [3] Bruce Xiaohan Nie, Caiming Xiong, and Song-Chun Zhu. Joint action recognition and pose estimation from video. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1293–1301, 2015.
- [4] Georgios Pavlakos, Xiaowei Zhou, Konstantinos G Derpanis, and Kostas Daniilidis. Coarse-to-fine volumetric prediction for single-image 3d human pose. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7025–7034, 2017.
- [5] Dario Pavllo, Christoph Feichtenhofer, David Grangier, and Michael Auli. 3d human pose estimation in video with temporal convolutions and semi-supervised training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7753–7762, 2019.
- [6] Heng Wang and Cordelia Schmid. Action recognition with improved trajectories. In *Proceedings of the IEEE international conference on computer vision*, pages 3551–3558, 2013.
- [7] Angela Yao, Juergen Gall, and Luc Van Gool. Coupled action recognition and pose estimation from multiple views. *International journal of computer vision*, 100(1):16–37, 2012.
- [8] Xiaojin Zhu and Andrew B Goldberg. Introduction to semi-supervised learning. *Synthesis lectures on artificial intelligence and machine learning*, 3(1):1–130, 2009.